

Black Box Variational Inference

Rajesh Ranganath, Sean Gerrish, David M. Blei

Presented by Dingcheng Hu

1

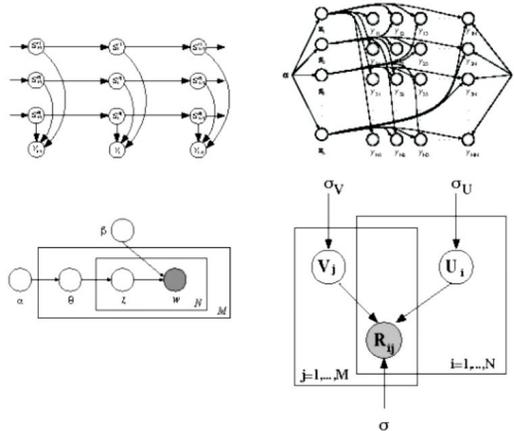
Problem

- Deriving a variational algorithm generally requires significant **model-specific** analysis
- These efforts hinder us from quickly developing and exploring a variety of models for a problem at hand.

2

Goal

Provide a **simple** way to approximate posteriors **quickly** in a **broad** class of latent variable models



Variational Inference

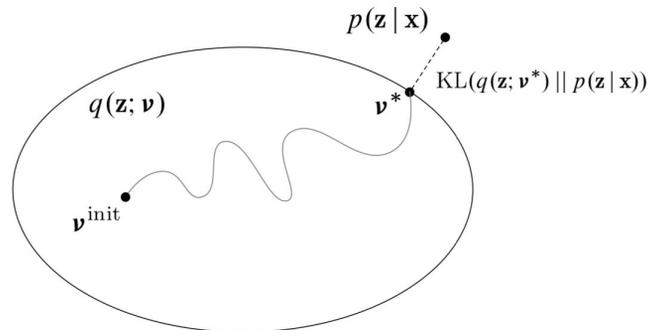
- A probabilistic model is a joint distribution of hidden variable z and observed variables x
- Inference about the unknowns is through the posterior, the conditional distribution of the hidden variables given the observations

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})}$$

- For most interesting models, the denominator is not tractable. We appeal to approximate the posterior.

Variational Inference

- Variational inference solves posterior computation with optimization
- Define a variational family of distributions over the latent variables $q(\mathbf{z}; \mathbf{v})$
- Update the parameter to fit the variational distribution to be close (in KL) to the exact posterior



5

Variational Inference

- Optimization is over a family q
- Find q that minimizes $\text{KL}(q || p(\mathbf{z} | \mathbf{x}))$
- But KL is intractable; VI optimizes the evidence lower bound (ELBO) instead:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_{\lambda}(z)} [\log p(\mathbf{x}, z) - \log q(z)]$$

- **The first term** rewards variational distributions that place high mass on configurations of the latent variable that also explain the observations
- **The second term** encourages variational distribution to be diffuse
- Minimizing the KL is the same as Maximizing the ELBO

6

Problem

- ELBO:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q_\lambda(z)}[\log p(x, z) - \log q(z)]$$

- Variational inference algorithms are normally derived by computing expectations and gradients
- Expectations often have no closed-form representation
- Computing the required expectations becomes intractable
- Require various expectation for each new model
- How do we get around this?

We want

Black Box Variational Inference

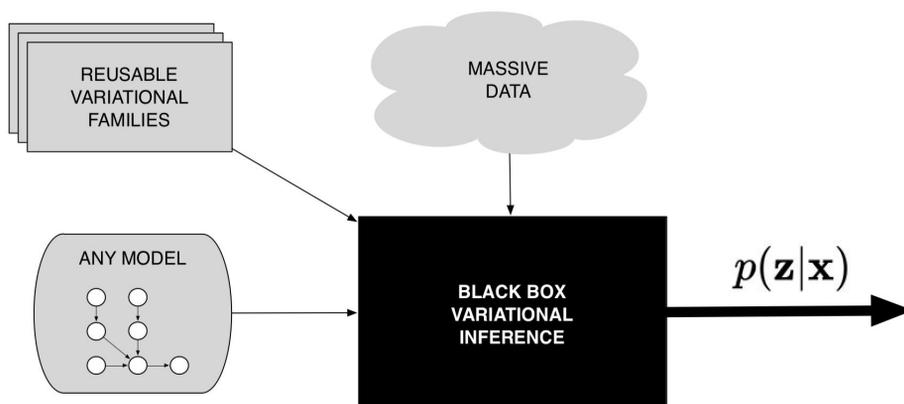
Why do we want Black Box Variational Inference

- Easily use variational inference with **any model**
- Minimal mathematical work beyond specifying the model $p(\mathbf{z}, \mathbf{x})$

9

General Idea

- Avoid computing model specific expectations and gradients
- Construct noisy gradients by sampling the variational family



10

Stochastic Optimization

Stochastic optimization maximizes a function using noisy gradients of that function. Formally:

- Objective: f
- Parameters: x
- A noisy gradient H ; $E[H] = \nabla f$
- Step size: ρ_t ;
- Update: $x_{t+1} \leftarrow x_t + \rho_t H(x_t)$

Converges to a local optimum when:

$$\begin{aligned}\sum_{t=1}^{\infty} \rho_t &= \infty \\ \sum_{t=1}^{\infty} \rho_t^2 &< \infty\end{aligned}$$

11

A Noisy Gradient of the ELBO

- The ELBO: $\mathcal{L}(\lambda) = \mathbb{E}_{q(z|\lambda)}[\log p(x, z) - \log q(z)]$
- Gradient of the ELBO:

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \mathbb{E}_q[\nabla_{\lambda} \log q(z|\lambda)(\log p(x, z) - \log q(z|\lambda))]$$

- How do we construct a noisy gradient of the ELBO?
- Using samples from q
- Noisy Gradient:

$$\frac{1}{S} \sum_{s=1}^S \nabla_{\lambda} \log q(z_s|\lambda)(\log p(x, z_s) - \log q(z_s|\lambda)),$$

where $z_s \sim q(z|\lambda)$

12

A Noisy Gradient of the ELBO

The noisy gradient:

$$\frac{1}{S} \sum_{s=1}^S \nabla_{\lambda} \log q(z_s|\lambda) (\log p(x, z_s) - \log q(z_s|\lambda)),$$

where $z_s \sim q(z|\lambda)$

To compute the noisy gradient of the ELBO we need

- Sampling from $q(z)$
- Evaluating $\nabla_{\lambda} \log q(z_s|\lambda)$
- Evaluating $\log p(x, z)$ and $\log q(z)$

13

A Noisy Gradient of the ELBO

The noisy gradient:

$$\frac{1}{S} \sum_{s=1}^S \nabla_{\lambda} \log q(z_s|\lambda) (\log p(x, z_s) - \log q(z_s|\lambda)),$$

where $z_s \sim q(z|\lambda)$

Fully black box

- Evaluating $\log p(x, z)$ is akin to defining the model
- The computations about q can be shared across models

14

Algorithm 1 Black Box Variational Inference

```
1: Input: data  $x$ , joint distribution  $p$ , variational family  $q$ .
2: Initialize  $\lambda_{1:n}$  randomly,  $t = 1$ .
3: repeat
4:   // Draw  $S$  samples from  $q$ 
5:   for  $s = 1$  to  $S$  do
6:      $z[s] \sim q$ 
7:   end for
8:    $\rho = t$ -th value of a Robbins Monro sequence
9:    $\frac{\hat{\partial} \mathcal{L}}{\partial \lambda} = \frac{1}{S} \sum_{s=1}^S \nabla_{\lambda} \log q(z[s]|\lambda) (\log p(x, z[s]) - \log q(z[s]|\lambda))$ 
10:   $\lambda = \lambda + \rho \frac{\hat{\partial} \mathcal{L}}{\partial \lambda}$ 
11:   $t = t + 1$ 
12: until change of held out predictive likelihood is less than  $\epsilon$ 
```

15

Variance Control

- Theoretically this is guaranteed to converge to local optima.
- But unfortunately the variance of the estimator can be very high
- We have to control the variance of the gradient
- Two techniques
 - Rao-Blackwellization
 - Control variates

16

Intuition

- Variance reduction methods work by replace the function whose expectation is being approximated by Monte Carlo with another function that has the same expectation but smaller variance
- To estimate $\mathbf{E}_q[f]$ via Monte Carlo
- We compute the empirical average of \hat{f}

$$\mathbf{E}_q[f] = \mathbf{E}_q[\hat{f}]$$
$$\mathbf{Var}_q[f] > \mathbf{Var}_q[\hat{f}]$$

17

Technique: Rao-Blackwellization

- Replace a random variable with its conditional expectation
- Try to estimate $\mathbb{E}[J(\mathbf{X}, \mathbf{Y})]$ with Monte Carlo
- Compute $\mathbb{E}[J(\mathbf{X}, \mathbf{Y})|\mathbf{X}]$ then estimate via Monte Carlo
- Proof: Define $\hat{J}(\mathbf{X}) = \mathbb{E}[J(\mathbf{X}, \mathbf{Y})|\mathbf{X}]$
- We have $\mathbf{E}[\hat{J}(\mathbf{X})] = \mathbf{E}[J(\mathbf{X}, \mathbf{Y})]$
- $\mathbf{Var}(\hat{J}(\mathbf{X})) = \mathbf{Var}(J(\mathbf{X}, \mathbf{Y})) - \mathbf{E}[(J(\mathbf{X}, \mathbf{Y}) - \hat{J}(\mathbf{X}))^2]$

Rao-Blackwellize the i th component in a Black box manner

$$\frac{1}{S} \sum_{s=1}^S \nabla_{\lambda_i} \log q_i(z_s|\lambda_i) (\log p_i(x, z_s) - \log q_i(z_s|\lambda_i)),$$

where $z_s \sim q_{(i)}(z|\lambda)$.

18

Technique: Control Variates

- Replace with $\hat{f}(z) \triangleq f(z) - a(h(z) - E[h(z)])$
- a is a scalar chosen to minimize the variance
- h is a function of our choice
- $E_q[\hat{f}] = E_q[f]$
- $\text{Var}(\hat{f}) = \text{Var}(f) + a^2 \text{Var}(h) - 2a \text{Cov}(f, h)$
- Good h have high correlation with f

But in **black box** variational inference

- Set h as $\nabla_{\lambda_i} \log q_i(z_s | \lambda_i)$
- Simply because $E_q[\nabla_{\lambda_i} \log q_i(z_s | \lambda_i)] = 0$ for any q

Maintains black box nature of the algorithm

19

Variance Reduction

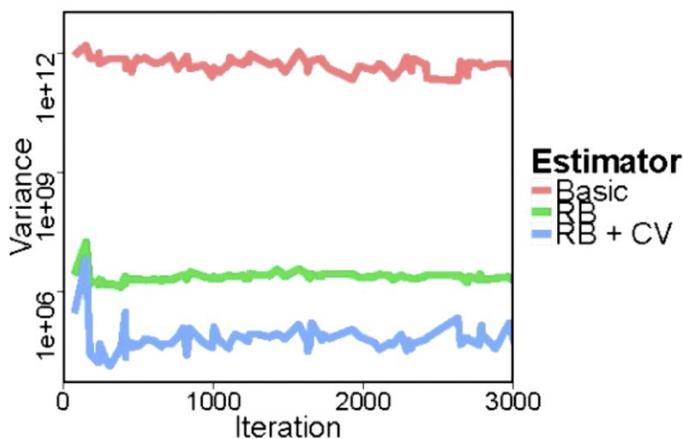


Figure: Variance comparison between our three estimators. The variance is reduced by several orders of magnitude.

20

Extension

- Set step size using AdaGrad
- Scalability by subsampling observations
 - Get a monte carlo estimate of $\log p(x, z)$

21

Experiment

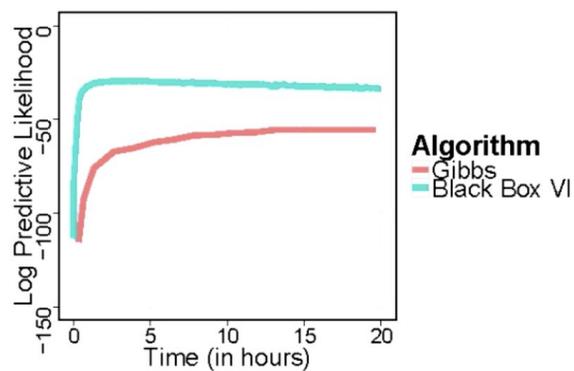


Figure: Comparison between Metropolis-Hastings within Gibbs and Black Box Variational Inference on a Gamma-Normal-TS model.

22

Summary

Black box variational inference only needs

- Log joint of the model $p(\mathbf{z}, \mathbf{x})$
- Computations of the variational approximation that can be shared across models